

Statistically optimal observables for global SMEFT fits

LHC EFT WG (Area 3)

24/10/23

Jaco ter Hoeve

VU Amsterdam & Nikhef theory group

in collaboration with R. G. Ambrosio,

M. Madigan, J. Rojo and V.Sanz

Where theory and experiment meet

We are progressively moving through the simulation chain (latent space)

$$p(x|c) \sim \int dz_{\text{det}} dz_{\text{shower}} dz_{\text{parton}} p(x|z_{\text{det}}) p(z_{\text{det}}|z_{\text{shower}}) p(z_{\text{shower}}|z_{\text{parton}}) p(z_{\text{parton}}|c)$$



SMEFT@NLO

(Likelihood free) Inference

Observed



Hidden/latent
variables



POI

Generation

Observed



Hidden/latent
variables



POI

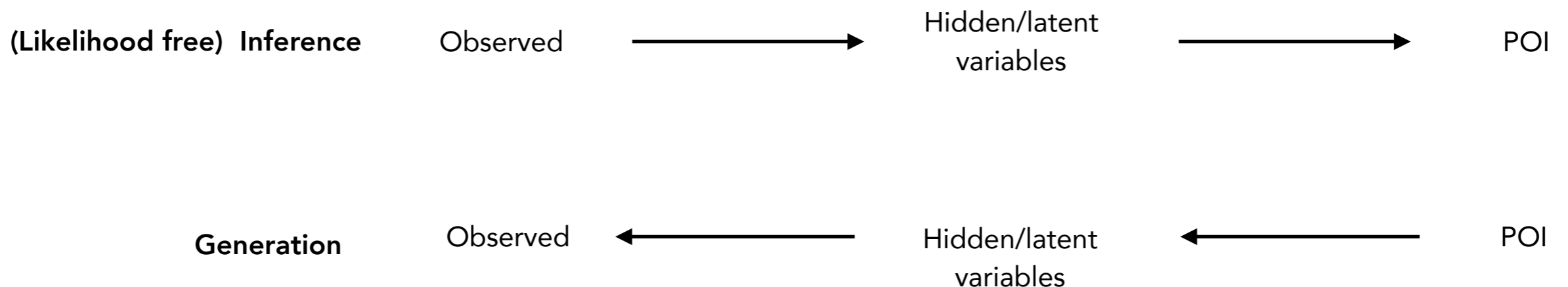
Where theory and experiment meet

We are progressively moving through the simulation chain (latent space)

$$p(x|c) \sim \int dz_{\text{det}} dz_{\text{shower}} dz_{\text{parton}} p(x|z_{\text{det}}) p(z_{\text{det}}|z_{\text{shower}}) p(z_{\text{shower}}|z_{\text{parton}}) p(z_{\text{parton}}|c)$$



Unbinned unfolding, Omnifold [1911.09107]



Likelihood free inference

- Starting from two balanced datasets \mathcal{D}_{SM} and \mathcal{D}_{EFT} drawn from $f(\mathbf{x} | \text{SM})$ and $f(\mathbf{x} | \text{EFT})$, we minimise e.g. the cross-entropy loss

$$L[g(\mathbf{x})] = -\frac{1}{N} \sum_{e \in \mathcal{D}_{\text{EFT}}} w_e \log(1 - g(\mathbf{x}_e)) - \frac{1}{N} \sum_{\mathcal{D}_{\text{SM}}} w_e \log g(\mathbf{x}_e)$$

Event weights
 $\{m_{\tilde{t}}, \eta_t, \Delta\phi, \dots\}$

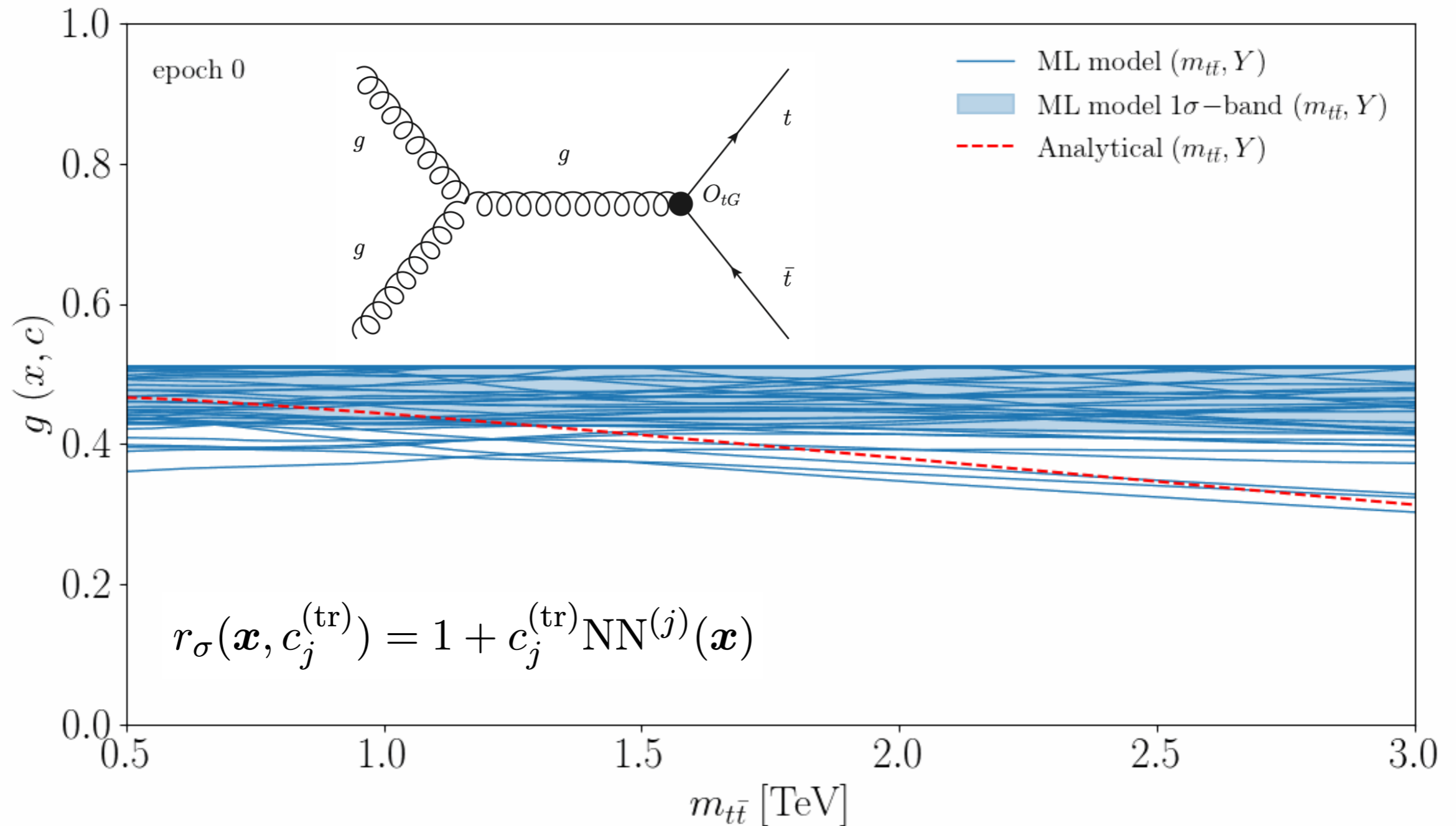
- The learned decision boundary $g(\mathbf{x})$ is one-to-one with the likelihood ratio (LR) as $N \rightarrow \infty$

$$\frac{\delta L}{\delta g} = 0 \implies \hat{g}(\mathbf{x}) = \left(1 + \frac{f(\mathbf{x} | \text{EFT})}{f(\mathbf{x} | \text{SM})} \right)^{-1} \equiv \frac{1}{1 + r(\mathbf{x})}$$

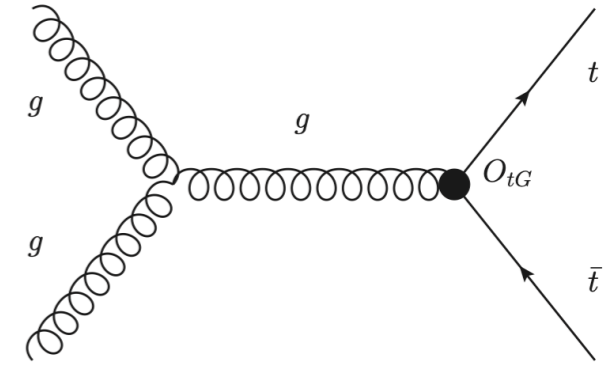
Parameterise with NNs

Likelihood free inference

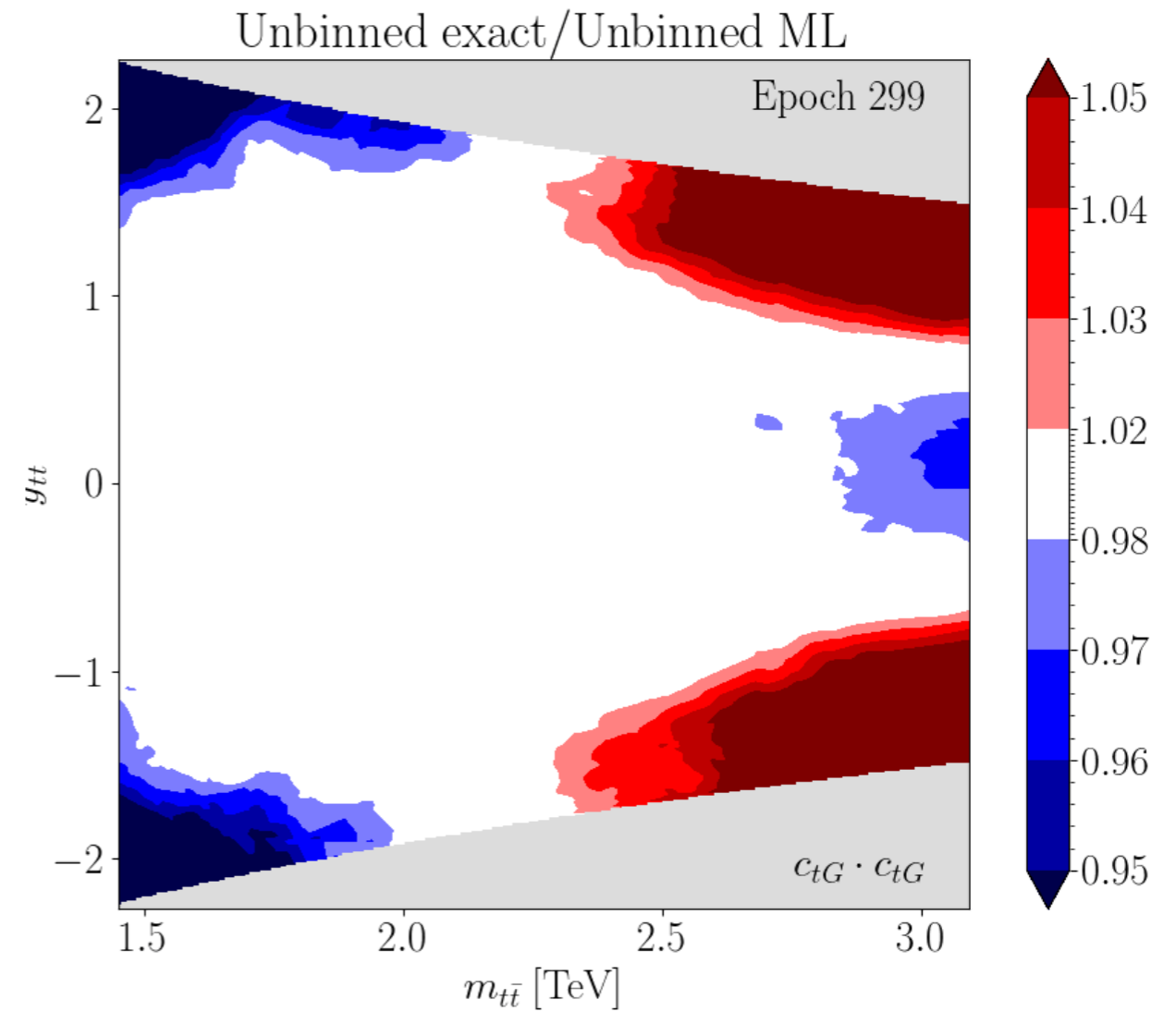
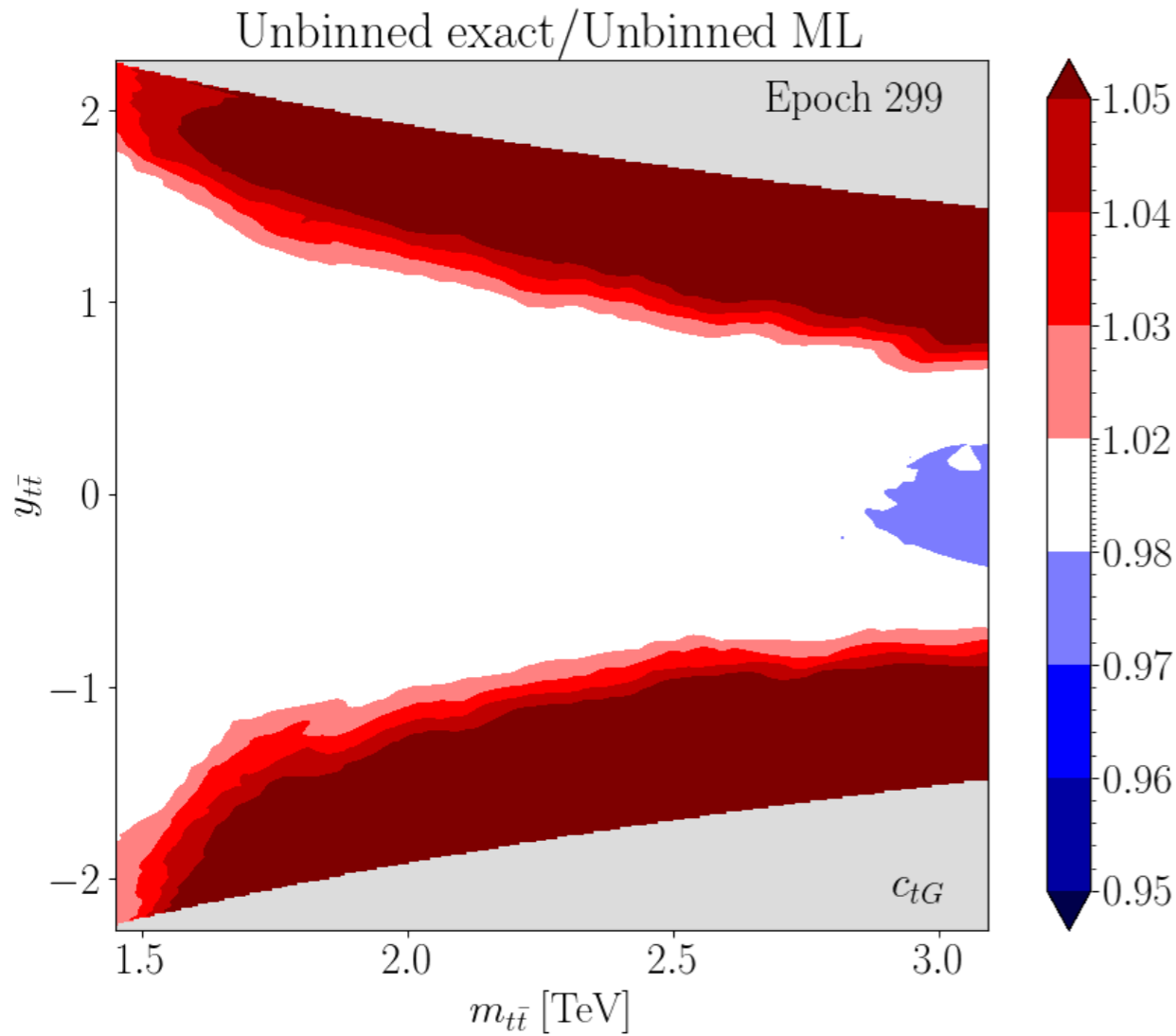
$$L[g(\mathbf{x}, \mathbf{c})] = -\frac{1}{N} \sum_{e \in \mathcal{D}_{\text{eff}}} w_e \log(1 - g(\mathbf{x}_e, \mathbf{c})) - \frac{1}{N} \sum_{e \in \mathcal{D}_{\text{sm}}} w_e \log g(\mathbf{x}_e, \mathbf{c})$$



Likelihood free inference



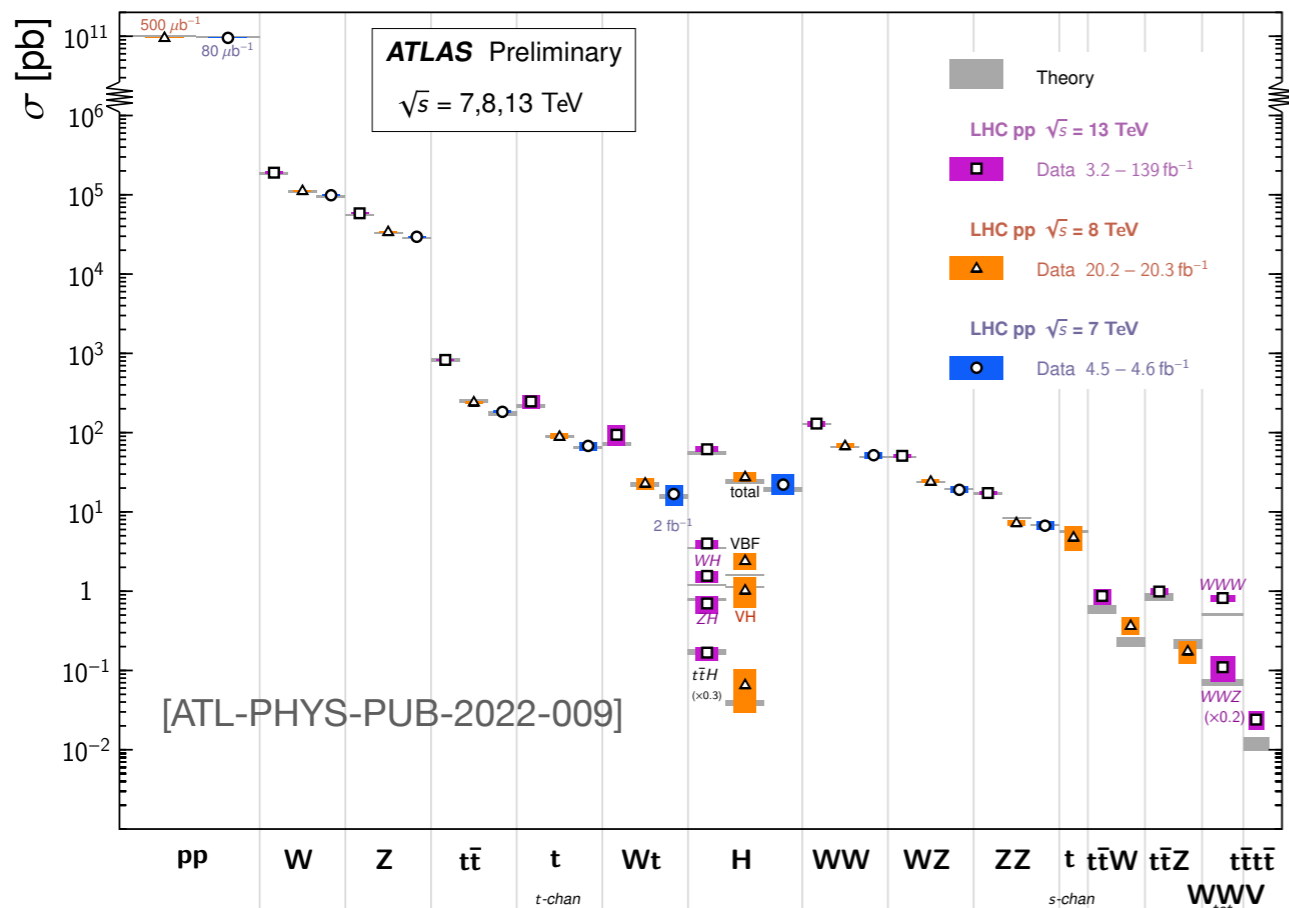
2 kinematic features



Why interesting for the SMEFT?

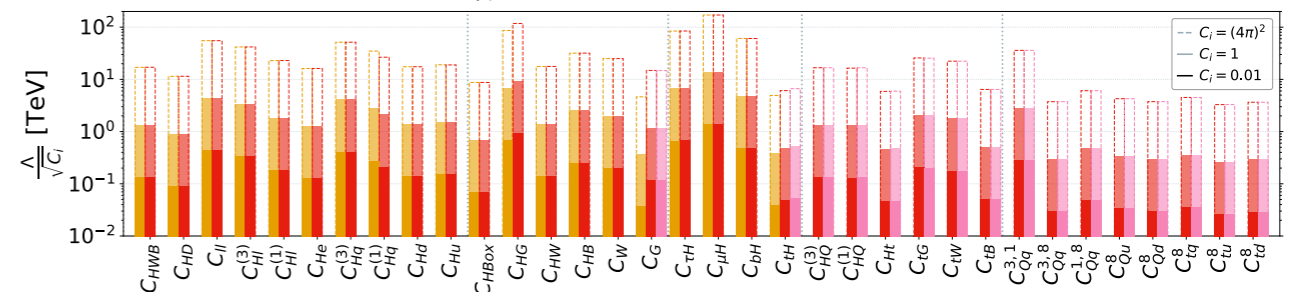
Standard Model Total Production Cross Section Measurements

Status: February 2022



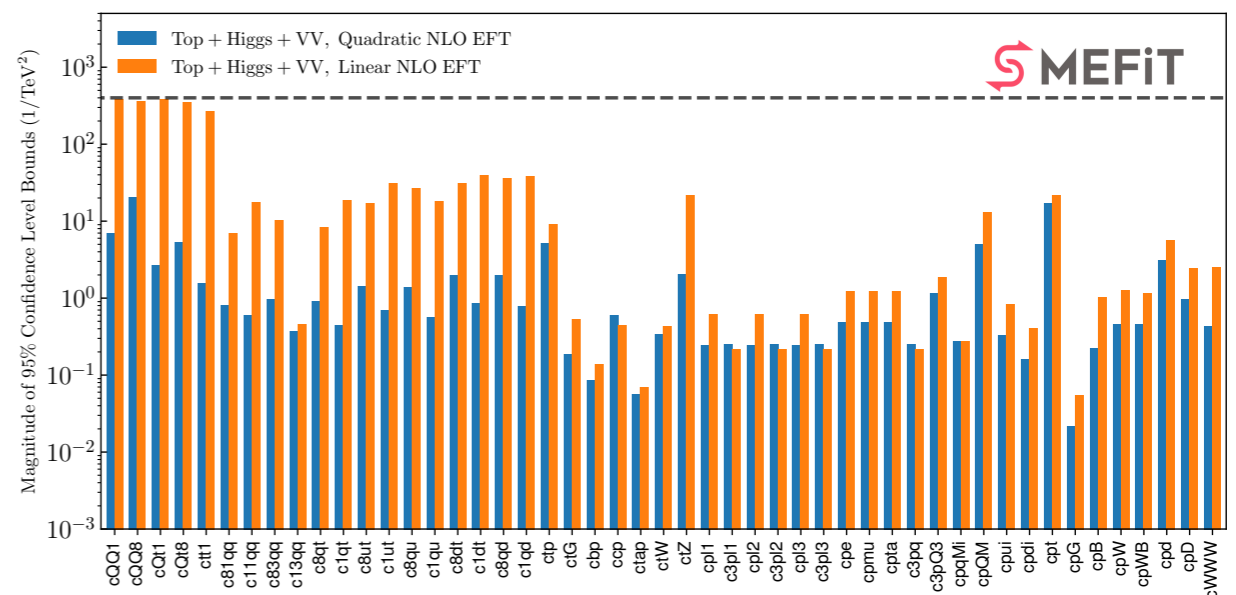
$$\chi^2 \sim \left(\sigma_i(c) - \sigma_{i,\text{exp}} \right) (\text{cov}^{-1})_{ij} \left(\sigma_j(c) - \sigma_{j,\text{exp}} \right)$$

[2012.02779]



[2302.0660]

$$\sigma(c) = \sigma_{\text{SM}} \left(1 + \sum_i^{N_{d6}} \kappa_i c_i + \sum_{i < j}^{N_{d6}} \tilde{\kappa}_{ij} c_i \cdot c_j \right)$$



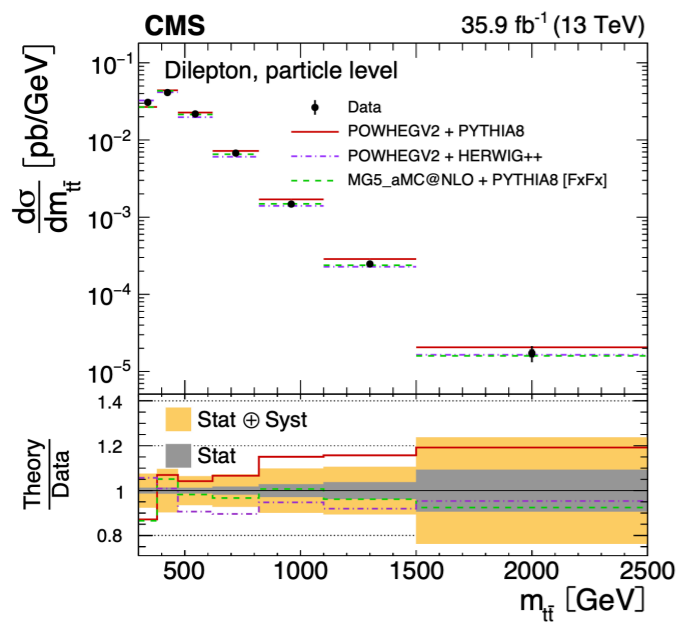
Why interesting for the SMEFT?

Standard Model Total Production Cross Section Measurements

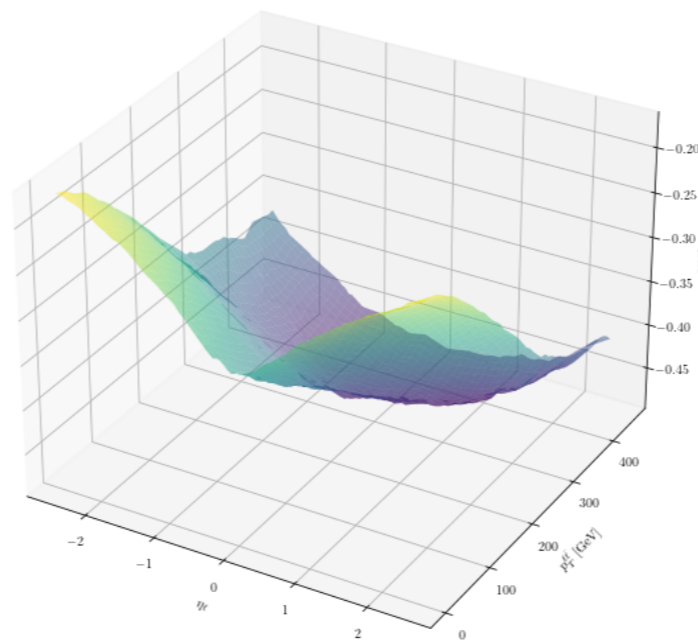
Status: February 2022



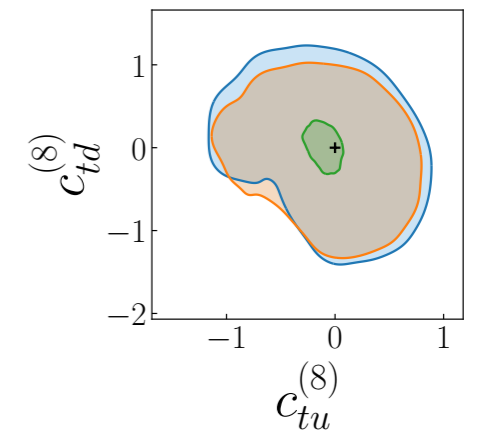
$$\chi^2 \sim \left(\sigma_i(c) - \sigma_{i,\text{exp}} \right) (\text{cov}^{-1})_{ij} \left(\sigma_j(c) - \sigma_{j,\text{exp}} \right)$$



Binned, univariate

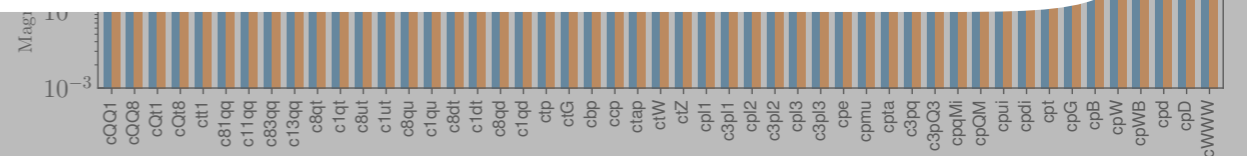


Unboxed, multivariate



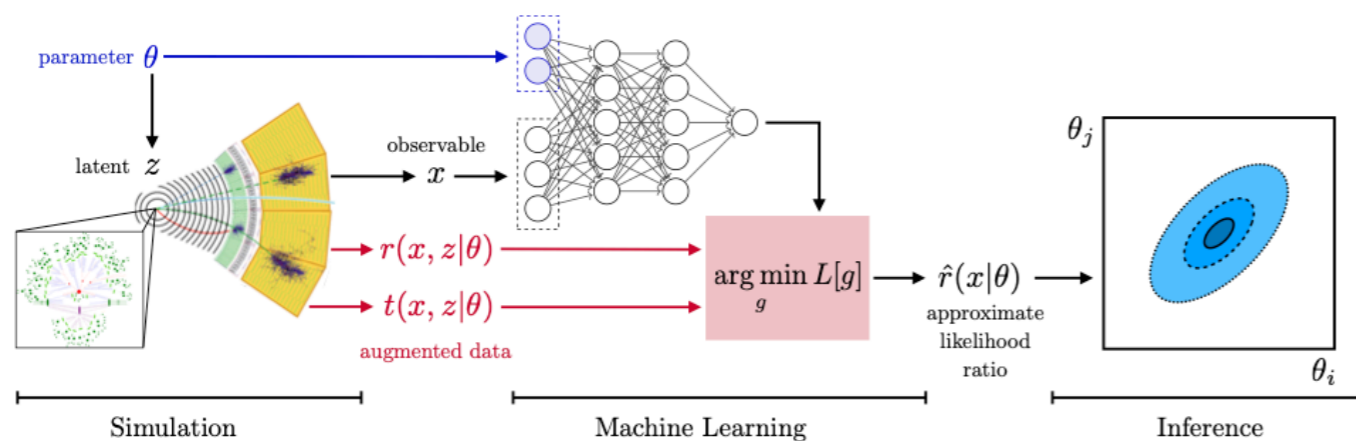
Optimal inference

$$\sigma(c) = \sigma_{\text{SM}} \left(1 + \sum_i K_i c_i + \sum_{i < j} K_{ij} c_i \cdot c_j \right)$$

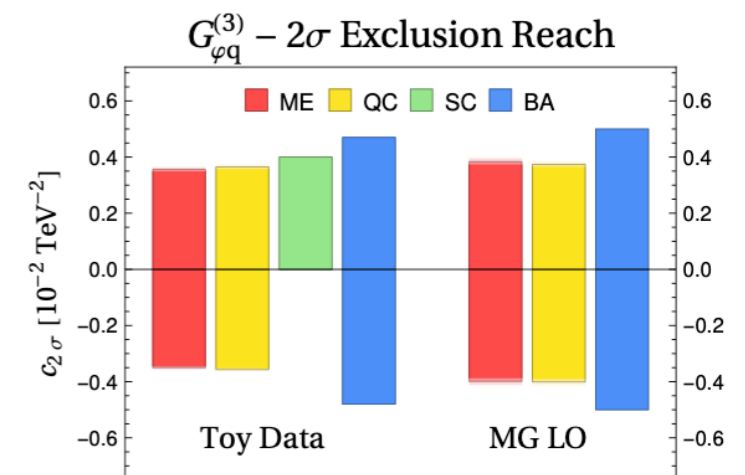


Related work

- ▶ Madminer series (J.Brehmer, K.Cranmer, G.Louppe et al.) [1907.10621, 1805.00020, ...]
- ▶ Parameterized classifiers for SMEFT (A. Glioti et al.) [2007.10356]
- ▶ Learning the EFT likelihood with tree boosting (R. Schöfbeck et al) [2205.12976]
- ▶ Back to the Formula (A. Butter, T. Plehn et al) [2109.10414]
- ▶ Boosted likelihood learning with event reweighing (A. Glioti et al) [2308.05704] See next talk!
- ▶ Designing Observables for Measurements with Deep Learning (O.Long, B. Nachman) [2310.08717]



[2010.06439]



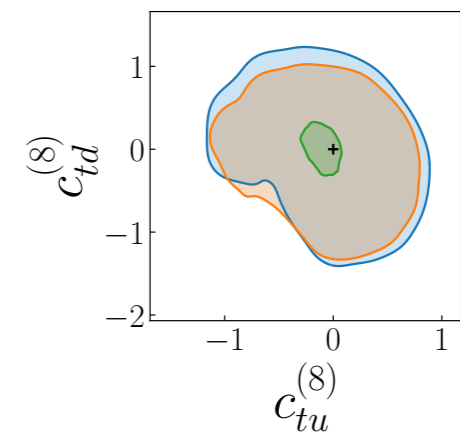
[2007.10356]

This talk

- Global efforts **reinterpret** existing “SM measurements” in an EFT context
- But which measurements are the most **sensitive** to EFT parameters?
 - Inclusive, single to multi-differential (which variables)
 - Binned or unbinned, which binning?

Framework needed to integrate unbinned multivariate observables into **global SMEFT fits**

- **Optimal** bounds on the EFT parameters
- Useful **diagnosis tool** to assess information loss



The ML4EFT framework

`pip install ml4eft`

<https://lhcfiteh.github.io/ML4EFT>

2211.02058 R. Gomez Ambrosio, JtH, M. Madigan, J. Rojo, V.Sanz

Open-source NN-based python framework for the integration of unbinned multivariate observables into global SMEFT fits

- ▶ **Goal:** provide optimal constraints on the SMEFT
- ▶ **Diagnostic tool:** what is the information loss incurred by a particular choice of bins?
- ▶ **Projections:** how will SMEFT constraints improve if unbinned data are made available?

The screenshot shows the ML4EFT documentation website. On the left is a navigation sidebar with a search bar and a tree view of the documentation structure. The main content area displays the documentation for the `ml4eft.core.classifier.Fitter` class. It includes the class signature, its base class (`object`), and its role as a training class. The `__init__` method is highlighted, showing its parameters: `json_path` (Path to json run card), `mc_run` (Replica number), `c_name` (EFT coefficient for which to learn the ratio function), `output_dir` (Path to where the models should be stored), and `print_log` (Set to true to print training progress to stdout, otherwise it prints to a log file only). Below the constructor, a list of methods is provided, including `load_data()`, `loss_fn(outputs, labels, w_e)`, `train_classifier(data_train, data_val)`, `training_loop(optimizer, train_loader, ...)`, and `weight_reset(m)`.

Modular structure, easy to maintain, well documented

Anticipating global fits

- Global EFT fits typically feature ~ 50 WCs and thus efficient scaling with the number of WCs becomes essential
- Solution: learn the coefficient functions separately and combine afterwards

$$r(\mathbf{x}, \mathbf{c}) \equiv \frac{d\sigma(\mathbf{x}, \mathbf{c})}{d\sigma(\mathbf{x}, \mathbf{0})} = 1 + \sum_{j=1}^{n_{\text{eft}}} r^{(j)}(\mathbf{x}) c_j + \sum_{j=1}^{n_{\text{eft}}} \sum_{k \geq j}^{n_{\text{eft}}} r^{(j,k)}(\mathbf{x}) c_j c_k$$

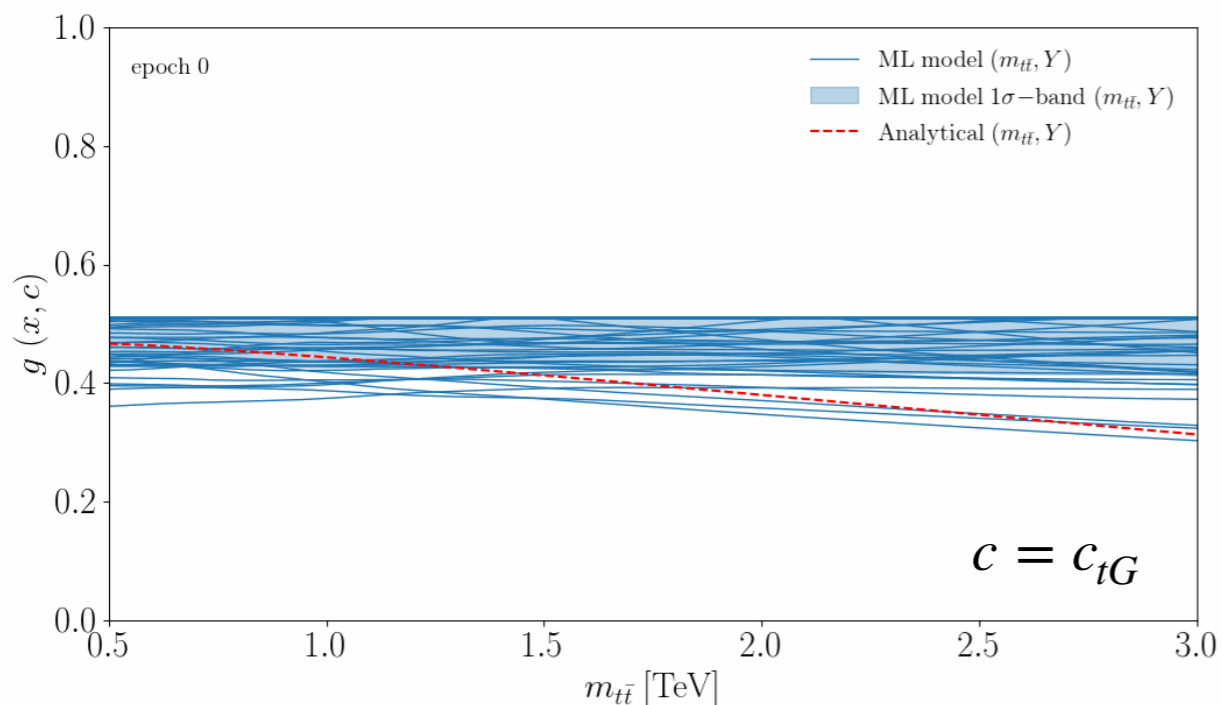
Example: to learn a single $r^{(j)}$, generate \mathcal{D}_{sm} and \mathcal{D}_{eft} at c_j up to $\mathcal{O}(\Lambda^{-2})$.
Then $r(\mathbf{x}, \mathbf{c}) = 1 + r^{(j)}(\mathbf{x}) c_j^{(\text{tr})}$ and training means

$$g(\mathbf{x}, c_j^{(\text{tr})}) = \left(1 + \left[1 + c_j^{(\text{tr})} \cdot \text{NN}^{(j)}(\mathbf{x}) \right] \right)^{-1} \quad \text{NN}^{(j)}(\mathbf{x}) \rightarrow r^{(j)}(\mathbf{x})$$

Uncertainty treatment

- ▶ Stick to a regime in which statistical uncertainties dominate over systematics
- ▶ Finite training data makes one subject to methodological uncertainties
- ▶ Solution: propagate uncertainties to the space of models by training multiple replicas

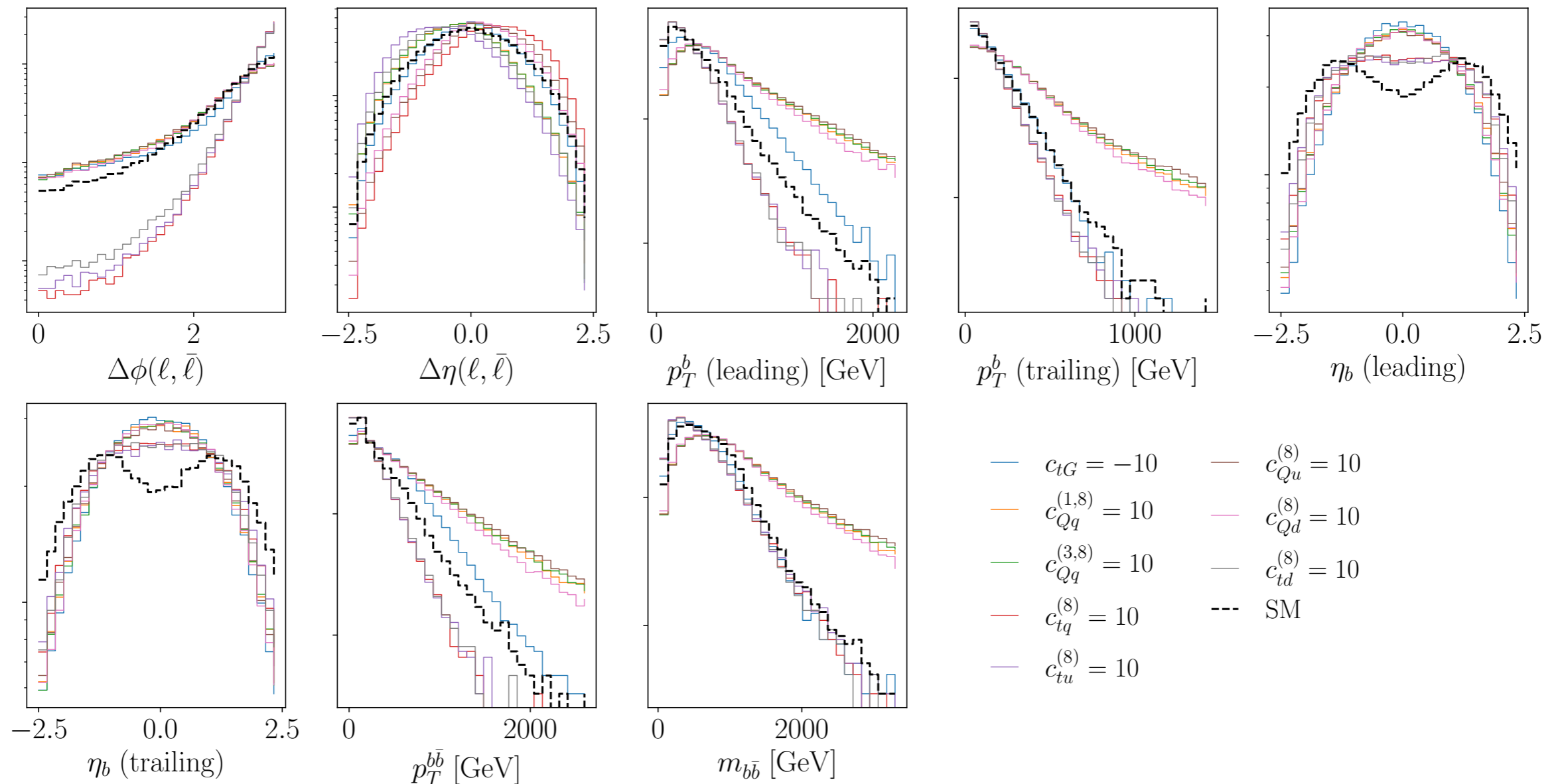
$$\hat{r}^{(i)}(\mathbf{x}, \mathbf{c}) \equiv 1 + \sum_{j=1}^{n_{\text{eff}}} \text{NN}_i^{(j)}(\mathbf{x})c_j + \sum_{j=1}^{n_{\text{eff}}} \sum_{k \geq j}^{n_{\text{eff}}} \text{NN}_i^{(j,k)}(\mathbf{x})c_j c_k, \quad i = 1, \dots, N_{\text{rep}}$$



Process	N_{rep}	\tilde{N}_{ev} (per replica)	N_{nn}	#trainings
$pp \rightarrow t\bar{t}$	50	10^5	4	200
$pp \rightarrow t\bar{t} \rightarrow b\bar{b}l^+\ell^-\nu_\ell\bar{\nu}_\ell$	25	10^5	40	1000

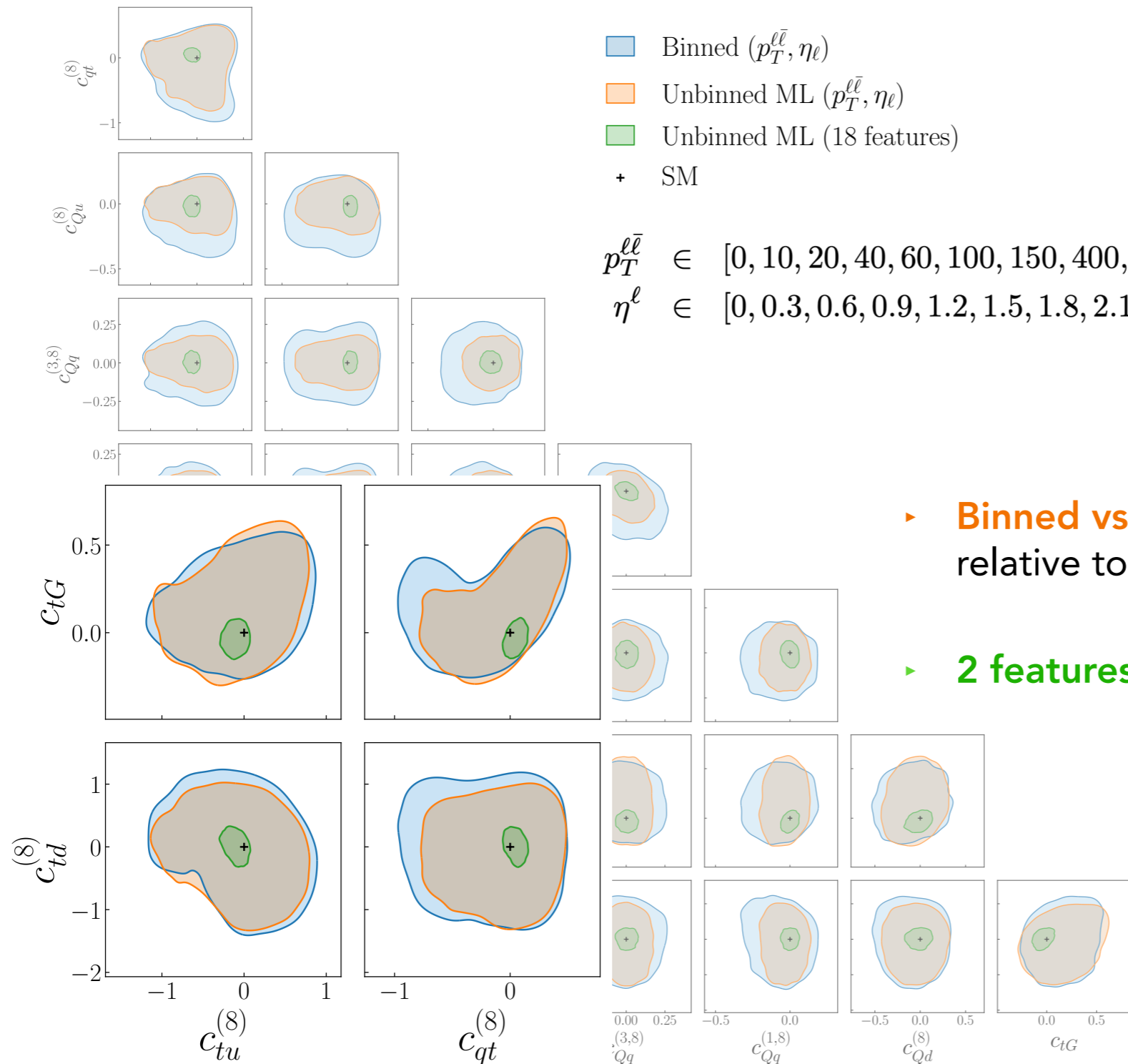
Let's go multivariate

- ▶ $pp \rightarrow t\bar{t} \rightarrow b\bar{b}\ell^+\ell^-\nu_\ell\bar{\nu}_\ell$: 18 features, 8 EFT coefficients
- ▶ $pp \rightarrow hZ \rightarrow b\bar{b}\ell^+\ell^-$: 7 features, 7 EFT coefficients



Unbinned observables in the top sector

Marginalised 95 % C.L. intervals, $\mathcal{O}(\Lambda^{-4})$ at $\mathcal{L} = 300 \text{ fb}^{-1}$

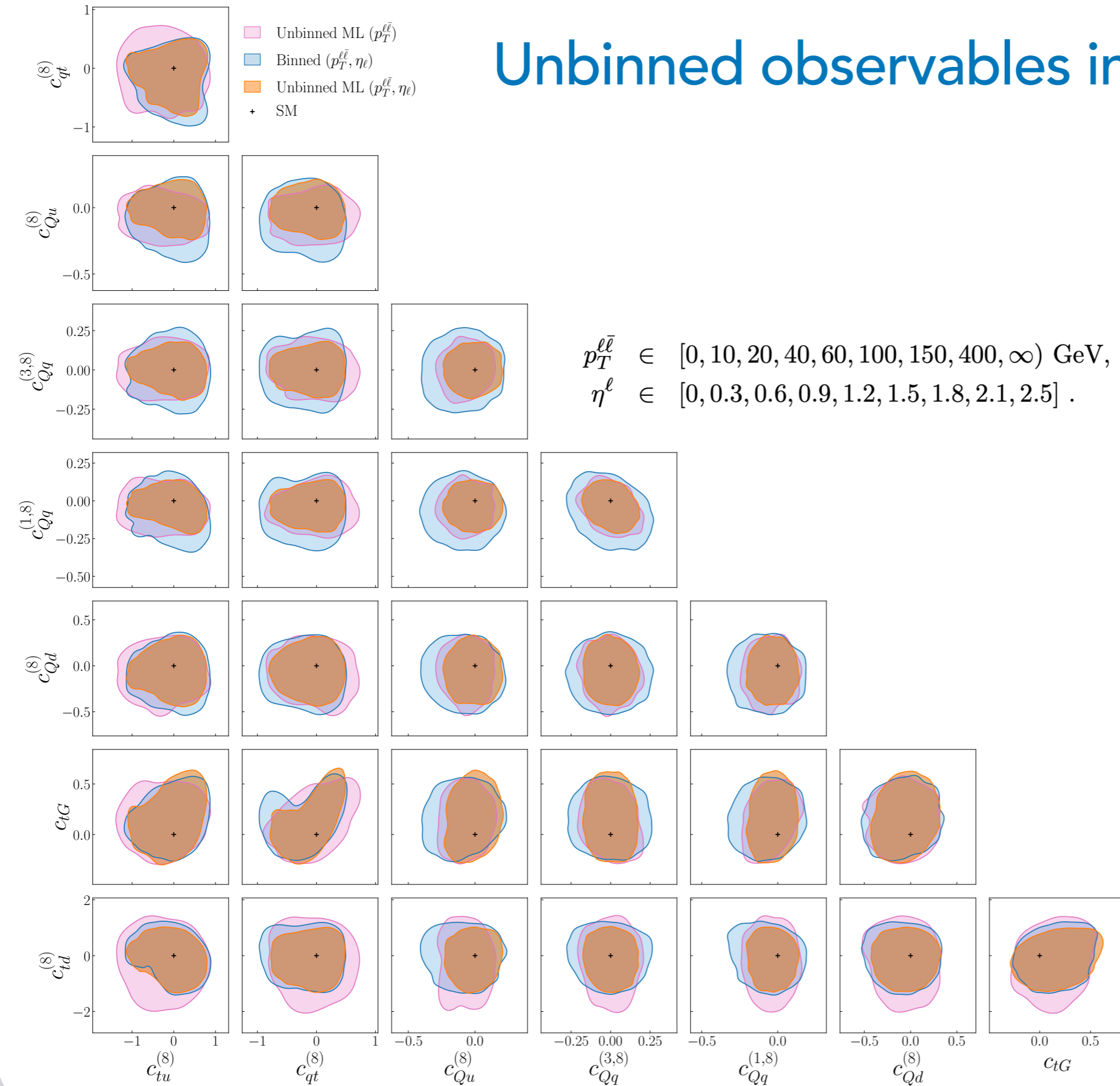


- ▶ **Binned vs unbinned** in $(p_T^{\ell\bar{\ell}}, \eta_\ell)$ small improvement relative to binned setup
- ▶ **2 features vs 18 features:** big increase in sensitivity

Unbinned observables in the top sector

$$pp \rightarrow t\bar{t} \rightarrow b\bar{b}\ell^+\ell^-\nu_\ell\bar{\nu}_\ell$$

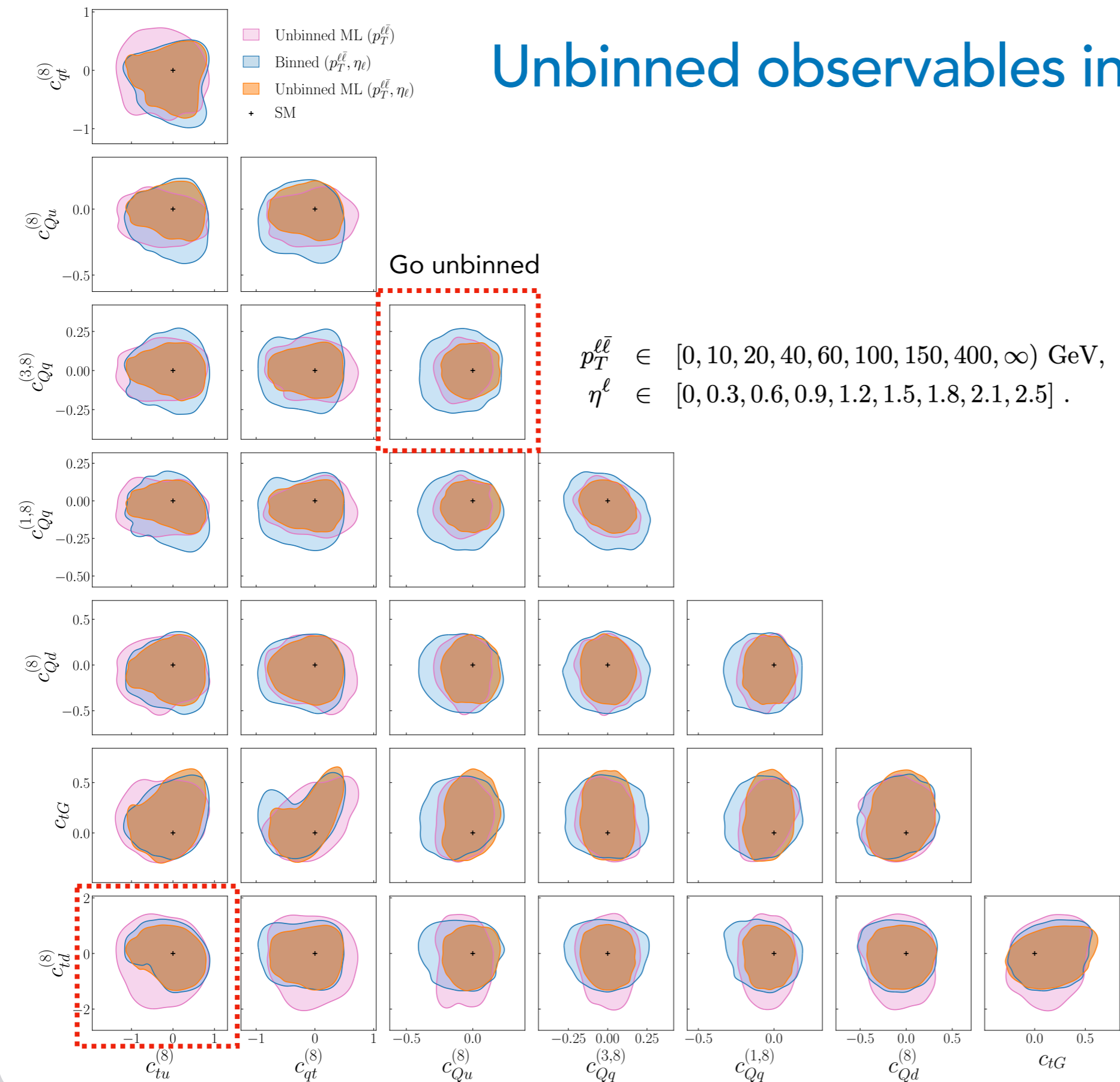
- ▶ Unbinned multivariate data is advantageous to constrain the EFT parameter space!
- ▶ Adding extra features or going unbinned?



Unbinned observables in the top sector

$$pp \rightarrow t\bar{t} \rightarrow b\bar{b}\ell^+\ell^-\nu_\ell\bar{\nu}_\ell$$

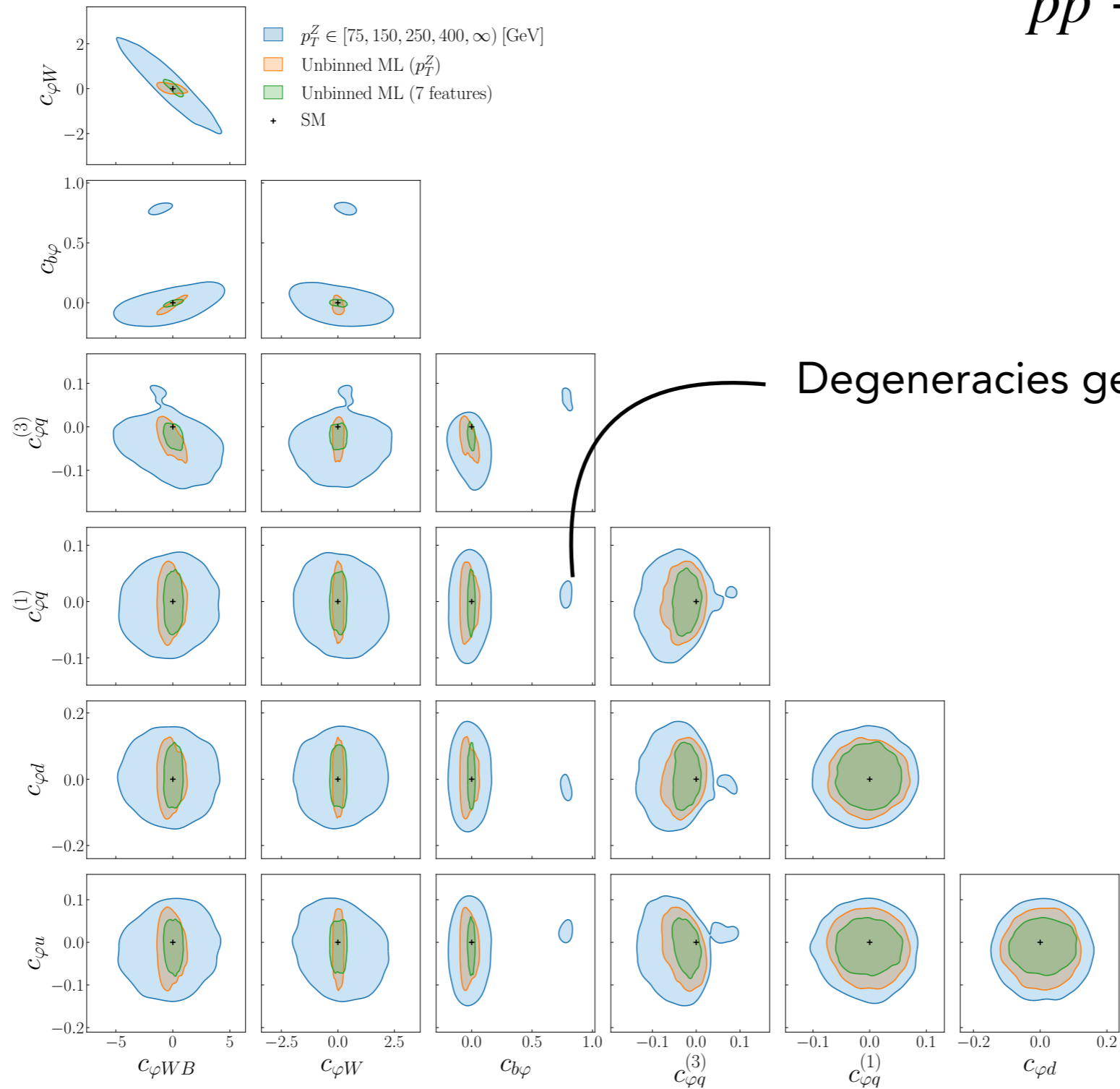
- ▶ Unbinned multivariate data is advantageous to constrain the EFT parameter space!
- ▶ Adding extra features or going unbinned?



Unbinned observables in Higgs + Z associated production

Marginalised 95 % C.L. intervals, $\mathcal{O}(\Lambda^{-4})$ at $\mathcal{L} = 300 \text{ fb}^{-1}$

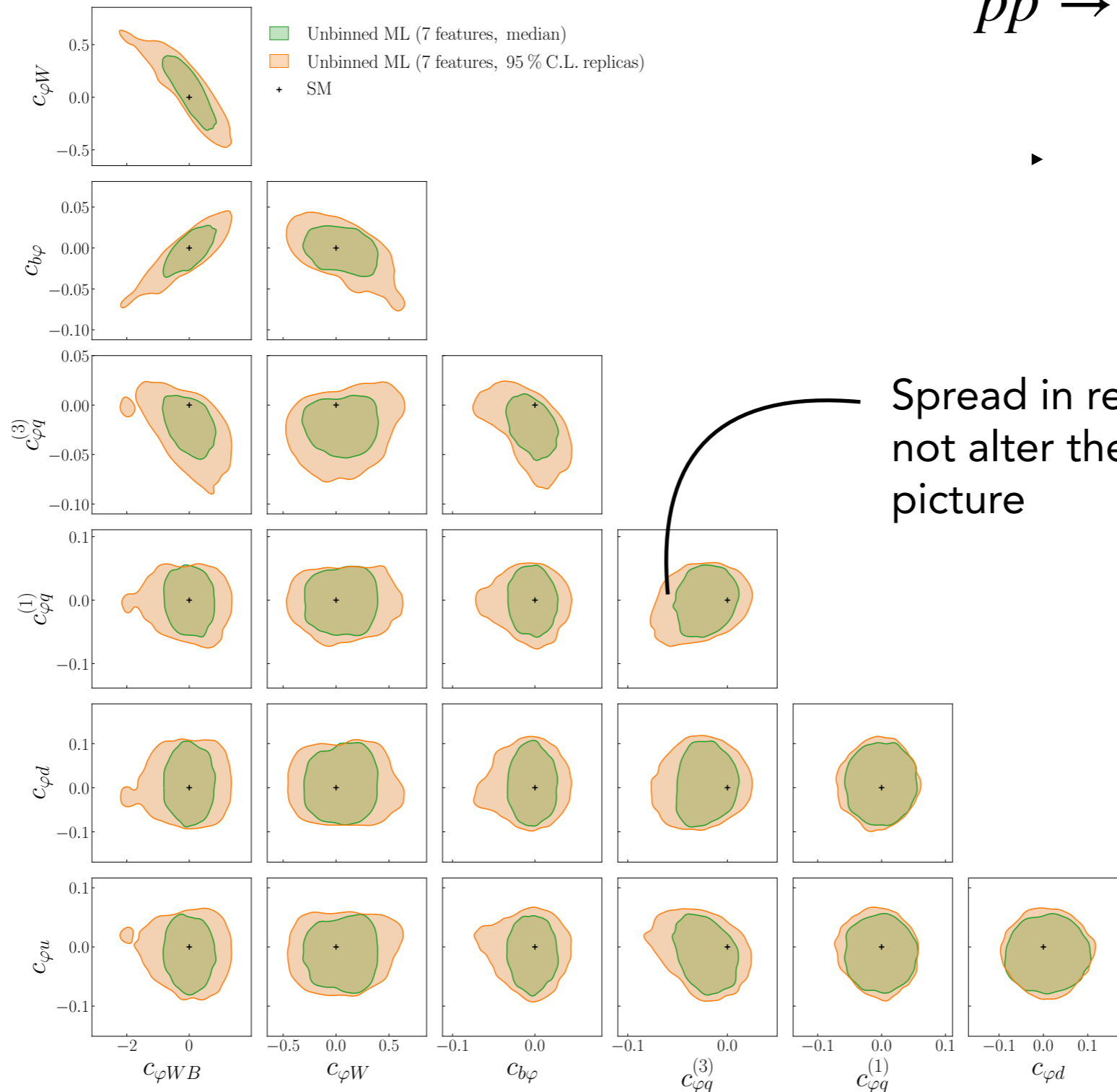
$$pp \rightarrow hZ \rightarrow b\bar{b}\ell^+\ell^-$$



Unbinned observables in Higgs + Z associated production

Marginalised 95 % C.L. intervals, $\mathcal{O}(\Lambda^{-4})$ at $\mathcal{L} = 300 \text{ fb}^{-1}$

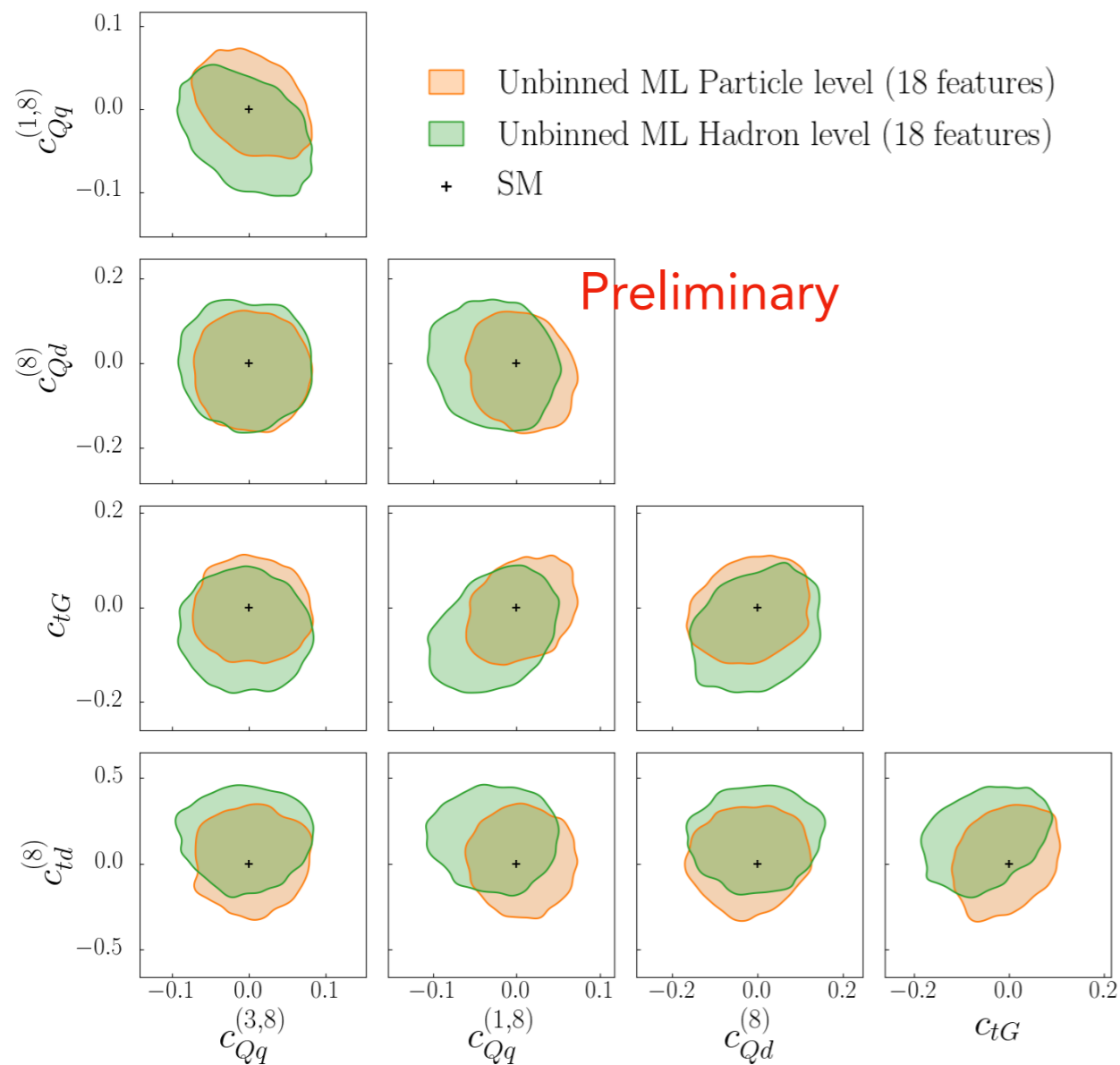
$$pp \rightarrow hZ \rightarrow b\bar{b}\ell^+\ell^-$$



Ongoing efforts

1. Hadronised level

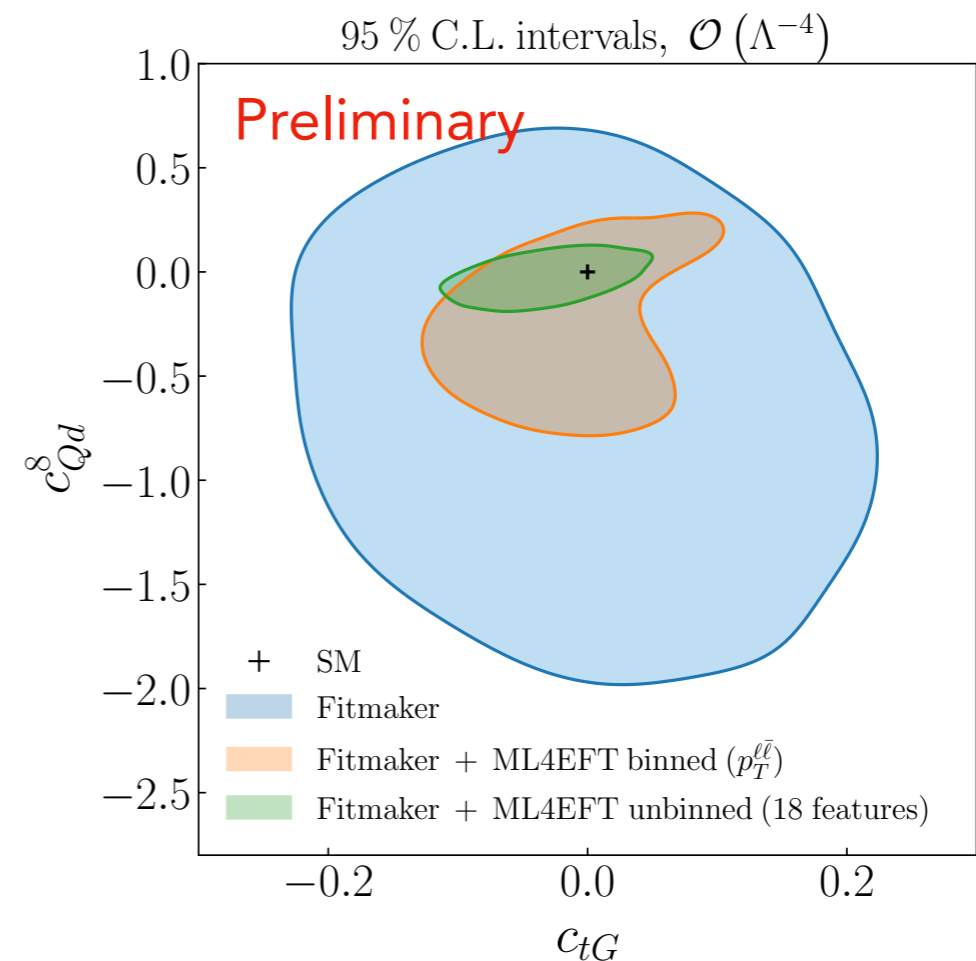
Marginalised 95 % C.L. intervals, $\mathcal{O}(\Lambda^{-4})$ at $\mathcal{L} = 300 \text{ fb}^{-1}$



MSc project by Pim Herbschleb

2. Integration into global fits

$$\log \mathcal{L}(c) = \sum_{k=1}^{N_D^{(\text{unbinned})}} \log \mathcal{L}_k^{\text{unbinned}}(c) + \sum_{k=1}^{N_D^{(\text{binned})}} \log \mathcal{L}_k^{\text{binned}}(c)$$



Conclusion and outlook

- Global EFT fits based on unbinned observables **enhance** the sensitivity significantly
- **ML4EFT** has efficient scaling properties, as required for global EFT fits and accounts for methodological uncertainties
- WIP: Integration into existing global fit frameworks
- WIP: benchmark study with A. Glioti et al.
- Please visit **ML4EFT** on GitHub (documentation + **tutorial**)

lhcfitefitnikhef.github.io/ML4EFT

Conclusion and outlook

- Global EFT fits based on unbinned observables **enhance** the sensitivity significantly
- **ML4EFT** has efficient scaling properties, as required for global EFT fits and accounts for methodological uncertainties
- WIP: Integration into existing global fit frameworks
- WIP: benchmark study with A. Glioti et al.
- Please visit **ML4EFT** on GitHub (documentation + **tutorial**)

lhcfitefitnikhef.github.io/ML4EFT

Thank you!